

# **Contents**

## 01. 프로젝트 개요

- 주제 및 선정 배경
- 주요 일정
- 스택 및 협업 툴
- 팀원 소개 및 역할

### 02. 개발 과정

- 데이터 정의
- 전처리
- 모델링
- 웹 개발

### 03. 결과

- 음악 추천 모델
- 웹 구현
- 시연 영상
- 한계점



## 01-1. 주제 및 선정 배경

### 최근 젊은 세대들 사이에서 취향, 상황, 분위기 등에 맞는 노래를 플레이리스트 형태로 소비하는 경향이 보임

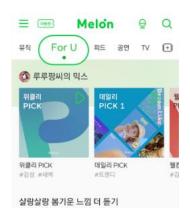


## 01-1. 주제 및 선정 배경

음악 스트리밍 어플 별로 다양한 방법을 도입하고 있으나, 여전히 곡 기반의 추천 시스템이 중심이 되고 있음

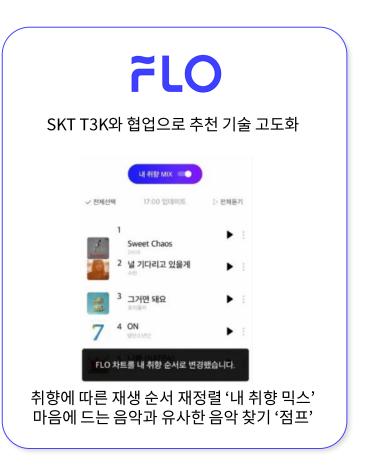
# Melon

방대한 데이터, 카카오 추천 엔진 도입 큐레이션에 T.P.O(시간, 장소, 상황) 반영



AI 기반 개인화 큐레이션 'For U' 나의 이용이력 분석 'MY'





## 01-1. 주제 및 선정 배경

플레이리스트 기반으로 음악을 추천해주는 '이곡어때' 시스템 개발을 주제로 선정함



# 01-2. 프로젝트 진행 일정

### 주요 업무 진행 일정은 아래와 같음

	8월 2주차	8월 3주차	8월 4주차	8월 5주차
아이디어 구상	8/10 8/11			
데이터 전처리	8/11	8/18		
모델 개발		8/16	8/24	
웹 구현			8/24 8/28	)
유지 보수			8/24 8/28	)
프로젝트 발표				8/29

# 01-3. 스택 및 협업 툴

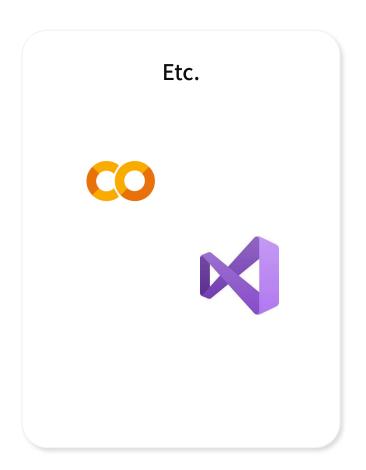
### 본 프로젝트를 위해 사용한 기술 및 협업 툴은 아래와 같음











# 01-4. 팀원 소개 및 역할

### 본 프로젝트는 아래 4인의 인원으로 구성됨



윤혜영 프로젝트 매니징 자연어 처리 및 모델 개발 웹 페이지 구현



조경희 전처리 및 모델 개발 웹 페이지 구현



**장윤식** 전처리 및 모델 개발 웹 서버 및 페이지구현



김 **빈** 웹 페이지 구현 모델 개발



# 02-1. 데이터 정의

### 카카오 아레나를 통해 다운 가능한 '멜론' 데이터를 이용함

### 플레이리스트 정보 데이터

플레이리스트 ID

플레이리스트 제목

태그 리스트

곡 리스트

좋아요 개수

수정 날짜

(115,071개 데이터)

### 곡 정보 데이터

곡 ID

앨범 ID

아티스트 ID 리스트

아티스트 리스트

곡 제목

곡 장르 리스트

곡 세부 장르 리스트

발매일

(707,989개데이터)

### 장르 데이터

장르 대분류

장르 소분류

# 02-2. 데이터 전처리

### 기존 데이터를 전처리한 결과는 아래와 같음

### 플레이리스트 정보 데이터

플레이리스트 ID

플레이리스트 제목

태그 리스트

곡 리스트

좋아요 개수

수정 날짜

(N) POS 플레이리스트명

(N) POS 태그

(N) 통합 태그

(N) 최종 태그

### 곡 정보 데이터

곡 ID

앨범 ID

아티스트 ID 리스트

아티스트 리스트

곡 제목

곡 장르 리스트

곡 세부 장르 리스트

발매일

(N) 곡 장르 태그

(N) 발매년도

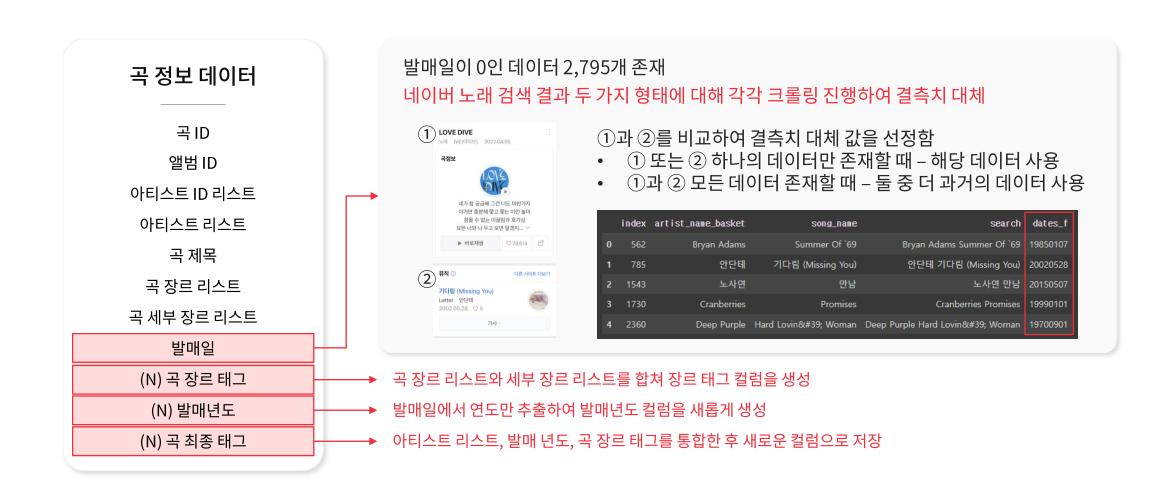
(N) 곡 최종 태그

#### 장르 데이터

장르 대분류

장르 소분류

## **02-2. 데이터 전처리** - 곡 정보 데이터



# 02-2. 데이터 전처리 - 플레이리스트 정보 데이터

### 플레이리스트 정보 데이터

플레이리스트 ID

플레이리스트 제목

태그 리스트

곡 리스트

좋아요 개수

수정 날짜

(N) POS 플레이리스트명

(N) POS 태그

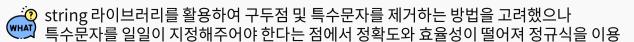
(N) 통합 태그

(N) 최종 태그

플레이리스트 제목이 같은 경우 비슷한 태그와 곡을 포함하고 있으므로 중복 처리하더라도 데이터 왜곡이 크지 않을 것이라 판단하여 제목 기준으로 중복처리

#### 정규식을 이용하여 플레이리스트 제목에서 특수문자를 제거함

train['plylst\_title'] = train['plylst\_title'].str.replace('[^가-힣A-Za-z0-9\&]','')



```
clean_columns = [] character_list = ['€', '£', '¥','。 • • □□', '□', '□□', '□', '□', '□', '花樣年華', '○', '□', '□'] import string for i in columns: tmp = i.translate(str.maketrans('', '', string.punctuation)).strip() for character in character_list: tmp = tmp.replace(character,'') clean_columns.append(tmp)
```

## 02-2. 데이터 전처리 - 플레이리스트 정보 데이터

### 플레이리스트 정보 데이터

플레이리스트 ID

플레이리스트 제목

태그 리스트

곡 리스트

좋아요 개수

수정 날짜

(N) POS 플레이리스트명

(N) POS 태그

(N) 통합 태그

(N) 최종 태그

### 특수문자를 제거한 플레이리스트 제목에서 형태소를 어근으로 추출

plylst\_tag = []
for i in range(len(train.plylst\_title)):
 tmp = okt.pos(train.plylst\_title[i], stem=True)
 plylst\_tag.append(tmp)

Komoran, Hannanum, Okt, Kkma를 이용해 같은 문장을 형태소 분석했을 때 Okt의 결과가 가장 적합하다고 판단하여 Okt를 활용하여 형태소 분석 및 어근 추출함

#### sentence = '○ 여행&드라이브 차안에서 함께들의면 기분좋아지는 리스트 sleep& drive○'

Komoran

Hannanum

Okt

Kkma

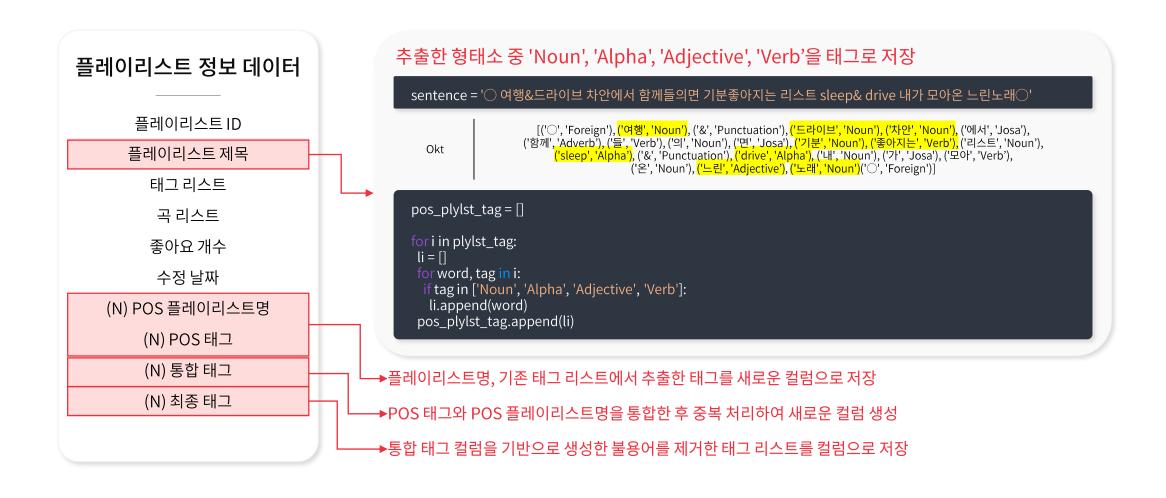
[('○', 'SW'), ('여행', 'NNP'), ('&', 'SW'), ('드라이브', 'NNP'), ('차', 'NNB'), ('안', 'NNP'), ('에서', 'JKB'), ('함께', 'MAG'), ('들', 'XSN'), ('의', 'JKG'), ('면', 'NNG'), ('기분', 'NNG'), ('좋', 'VA'), ('아', 'EC'), ('지', 'VX'), ('는', 'ETM'), ('리스트', 'NNP'), ('sleep', 'SL'), ('&', 'SW'), ('drive', 'SL'), ('○', 'SW')]

[('○', 'N'), ('여행&드라이브', 'N'), ('차안', 'N'), ('에서', 'J'), ('함께들의', 'N'), ('이', 'J'), ('면', 'E'), ('기분좋', 'P'), ('아', 'E'), ('지', 'P'), ('는', 'E'), ('리스트', 'N'), ('sleep&drive○', 'N')]

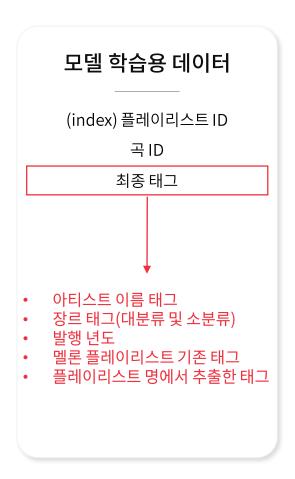
[('○', 'Foreign'), ('여행', 'Noun'), ('&', 'Punctuation'), ('드라이브', 'Noun'), ('차안', 'Noun'), ('에서', 'Josa'), ('함께', 'Adverb'), ('들', 'Verb'), ('의', 'Noun'), ('면', 'Josa'), <mark>('기분', 'Noun'), ('좋아지는', 'Verb'),</mark> ('리스트', 'Noun'), ('sleep', 'Alpha'), ('&', 'Punctuation'), ('drive', 'Alpha'), ('○', 'Foreign')]

[('○', 'SW'), ('여행', 'NNG'), ('&', 'SW'), ('드라이브', 'NNG'), ('차안', 'NNG'), ('에서', 'JKM'), ('함께', 'MAG'), ('들', 'VV'), ('ㄹ', 'ETD'), ('의', 'NNG'), ('면', 'NNG'), ('기분', 'NNG'), ('좋아지', 'VV'), ('는', 'ETD'), ('리스트', 'NNG'), ('sleep', 'OL'), ('&', 'SW'), ('drive', 'OL'), ('○', 'SW')]

# 02-2. 데이터 전처리 - 플레이리스트 정보 데이터

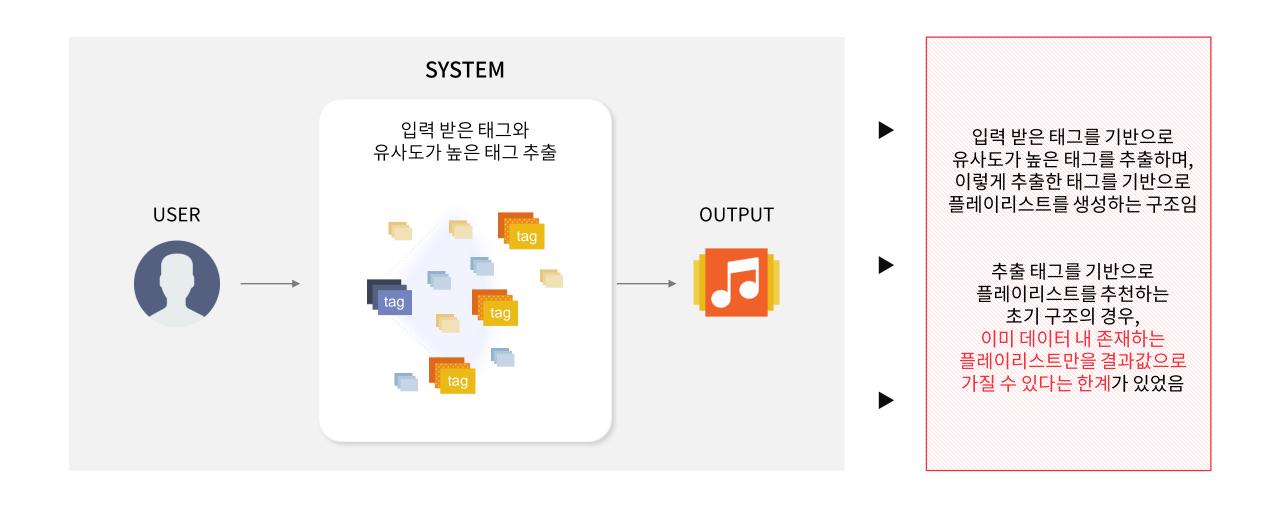


# 02-2. 데이터 전처리 – 모델 학습용 데이터

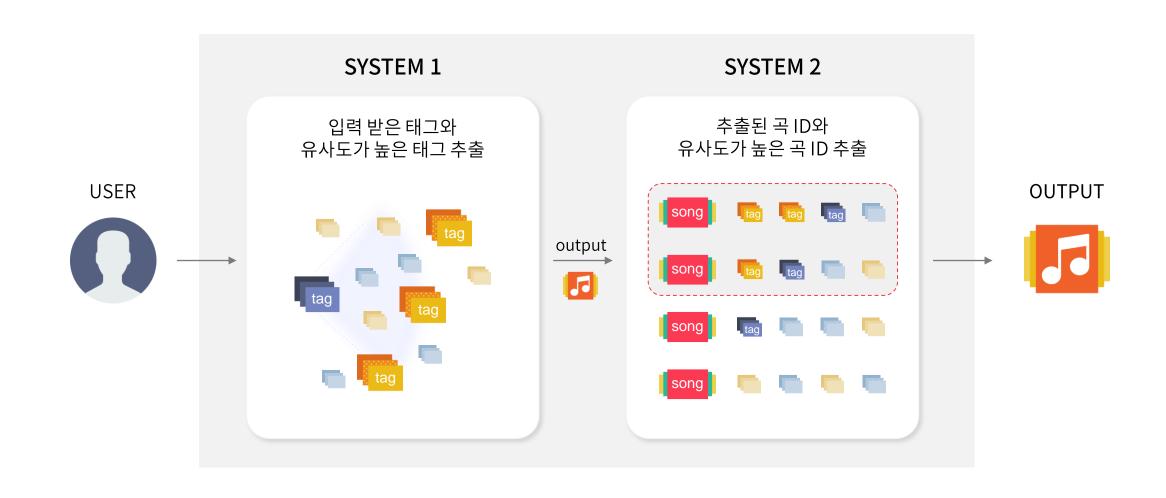


	song_id	final_tags			
0	0	[죽다, 회상, 아하, 2000년대, 영화음악, POP, 드라이브, 하루, 비, 추억]			
1	3	[느낌, 매력, 통해, 매장, 편집, 스트레스, 슬픔, 취향, 페스티벌, 성대, 목			
2	4	[2000년대, Jude Law, 이지리스닝, newage, 뉴에이지]			
3	5	[아름답다, 2000년대, 리메이크, 비, 메콩, 살, 기억, 페스티벌, 추억, 아			
4	6	[2000년대, 산책, San Francisco Symphony, 교향/관현악, 조			
615133	707984	[right, 90년대, 월드뮤직, nowglobal, Fela Kuti, pitc			
615134	707985	[80년대, POP, Cyndi Lauper, 로퍼, 라이벌, 마돈나, floor,			
615135	707986	[김광석, 2000년대, 매력, 들리다, 발라드, 리메이크, 업데이트, 이별, 설레			
615136	707987	[덥다, 새벽, 2000년대, 집중, 컴필레이션, 오후, 뉴에이지, Nature P			
615137	707988	['90, 록/메탈, 90년대, 샤우팅, 데뷔, 김경호, 지존]			
615138 rows × 2 columns					

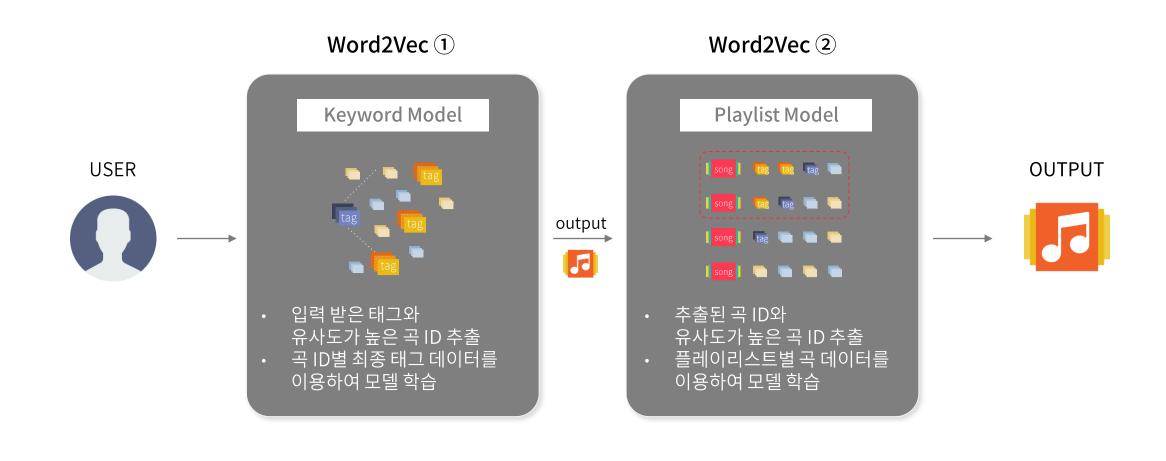
### 이곡어때 음악 추천 시스템 모델의 초기 구조는 다음과 같음



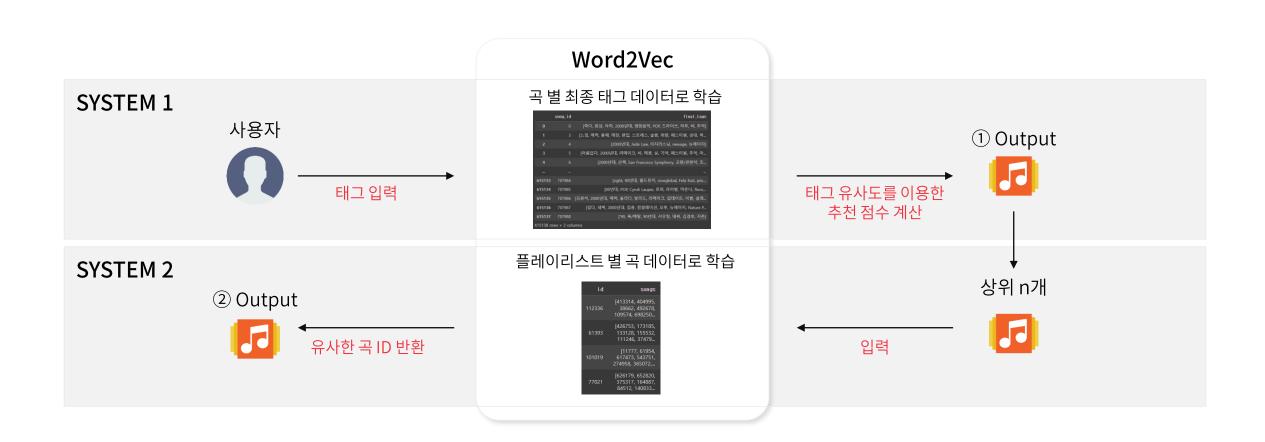
### 초기 구조를 보완하여 고안한 음악 추천 시스템 모델의 기본 구조는 다음과 같음



### 각 시스템별로 다른 데이터를 학습시킨 Word2Vec 모델을 결합하여 결과물을 도출함



### 이곡어때 음악 추천 시스템 모델의 전체 구조는 아래와 같음



### Word2Vec 모델과 FastText 모델을 비교하여 결과가 더 적합한 Word2Vec 모델을 선택함



FastText의 경우, 입력한 단어(예: 여름)가 포함된 단어의 유사도가 높게 나타나고 있는 반면, Word2Vec은 입력 단어가 포함되지 않더라도 의미가 유사한 경우 유사도가 높게 나타나는 것을 확인함

### Word2Vec 모델과 FastText 모델을 비교하여 결과가 더 적합한 Word2Vec 모델을 선택함

```
def exclusion(word_list):
    count = sum([score for _, score in word_list])
    avg = count / len(word_list)

    up_word = []
    for w, s in word_list: # w : 단어, s : 유사도
        if s > avg:
            up_word.append((w, s))
    return up_word

light = word2vect
light_word = exclight_word = exclight_word

[('밤메들기좋은도 ('깊은밤', 0.674 ('늦은밤', 0.674 ('늦은밤', 0.674 ('늦은밤', 0.674 ('늦은밤', 0.674 ('맞은'), 0.684 ('라고기', 0.6130 ('오지', 0.6130 ('오지', 0.6089 ('생각', 0.6089 ('생각', 0.6089 ('생각', 0.6089 ('고요한', 0.
```

```
light = word2vec.wv.most_similar('새벽', topn=20)
light_word = exclusion(light)
light_word

[('밤메들기좋은노래', 0.7161941528320312),
('깊은밤', 0.6749042272567749),
('늦은밤', 0.6706857681274414),
('밤공기', 0.6380821466445923),
('먹서사', 0.6227041482925415),
('밤새벽', 0.6182688474655151),
('오지', 0.6130034923553467),
('불면', 0.6089522242546082),
('생각', 0.6083134412765503),
('고요한', 0.6071279644966125)]
```

Word2Vec

모델 학습용 데이터

```
light = fasttext.wv.most_similar('새벽', topn=20)
light_word = exclusion(light)
light_word

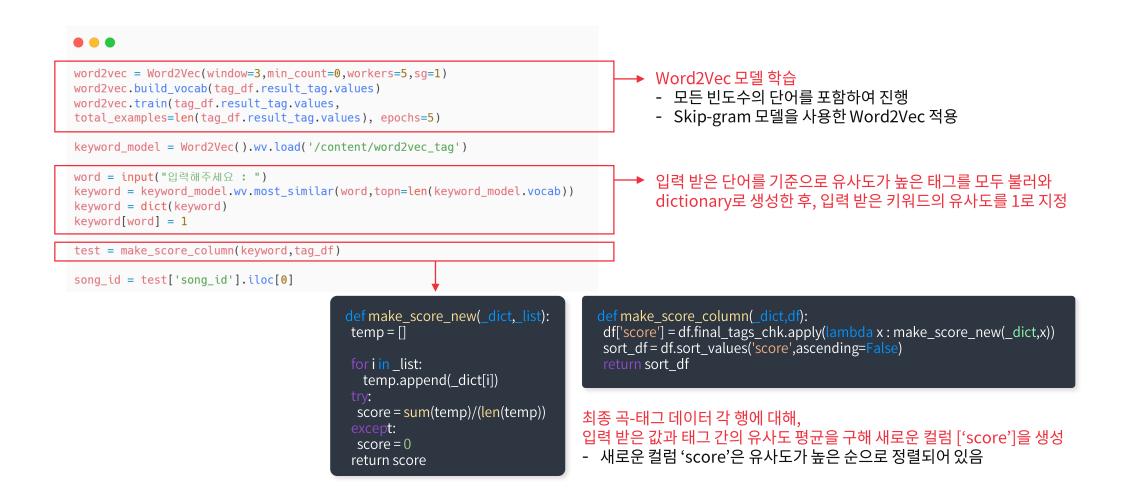
[('새벽새벽', 0.9829696416854858),
('깊은새벽', 0.9727602005004883),
('깊은새벽', 0.9721930623054504),
('새벽비', 0.9577507972717285),
('샘-새벽', 0.9570854306221008),
('새벽발', 0.9504075646400452),
('새벽플', 0.9479882121086121),
('새벽감', 0.9471359252929688)]
```

FastText

평균 이상의 유사도를 갖는 태그만 가져오는 exclusion 함수를 정의한 후, 상위 20개만 확인 → 이곡어때 음악 추천 시스템에는 Word2Vec을 채택하는 것이 적합하다고 판단하여 해당 모델을 선택함

# **02-3. 모델링** – Keyword Model

### 음악 추천 시스템 모델링을 위해 작성한 코드는 아래와 같음



# **02-3. 모델링** – Playlist Model

### 음악 추천 시스템 모델링을 위해 작성한 코드는 아래와 같음



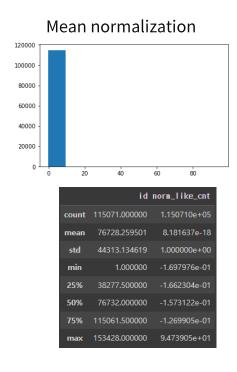


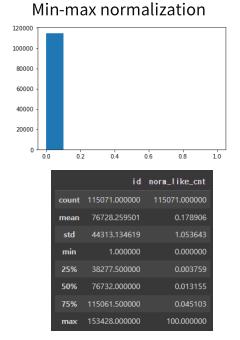
### 모델 최적화를 위해 score 모델을 고안하고자 함

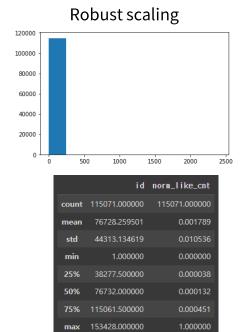
mean normalization, min-max normalization, robust scaling 등을 시도했지만, 데이터 편차를 줄이지 못함

#### 모델 학습용 데이터 편차 매우 큼

	id	like_cnt
count	115071.000000	115071.000000
mean	76728.259501	95.197687
std	44313.134619	560.653757
min	1.000000	0.000000
25%	38277.500000	2.000000
50%	76732.000000	7.000000
75%	115061.500000	24.000000
max	153428.000000	53211.000000







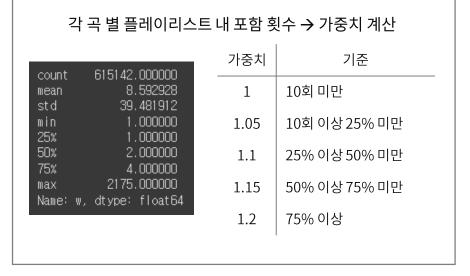
### 모델 최적화를 위해 score 모델을 고안하고자 함

이에 플레이리스트 별 좋아요 개수, 각 곡 별 플레이리스트 내 포함 횟수의 분포를 고려하여 임의로 score model을 개발

플레이리스트 별 좋아요 개수를 기준으로 곡에 점수 부여 → 평균 점수 계산

	id	like_cnt	점수	기준
count	115071.000000	115071.000000	1	250/ 1111
mean	76728.259501	95.197687	1	25% 미만
std	44313.134619	560.653757	2	25% 이상 50% 미만
min	1.000000	0.000000	_	
25%	38277.500000	2.000000	3	50% 이상 75% 미만
50%	76732.000000	7.000000	4	75% 이상 mean값 미만
75%	115061.500000	24.000000	·	1979   0 111641142   1
max	153428.000000	53211.000000	5	mean값 이상
			'	





최종 가중치를 기준으로 데이터를 정렬했을 때, 특정 가수의 노래가 지나치게 높은 가중치를 가지는 것을 확인할 수 있었음 → 해당 가중치를 모델에 적용하면 데이터 왜곡으로 이어질 것을 우려하여 기존 모델을 사용하는 것으로 결정함

## 02-4. 웹 개발

### 이곡어때 음악 추천 시스템을 웹으로 구현할 때 Django를 이용함



- MVT(Model-View-Template) 구조
- 파이썬을 서버측 언어로 사용할 수 있음
- 사전 생성된 파이썬 파일 구조를 이용할 수 있음
- ORM(Object-Relational Mapper)을 내장하고 있어 다양한 데이터베이스에 유연하게 액세스 가능

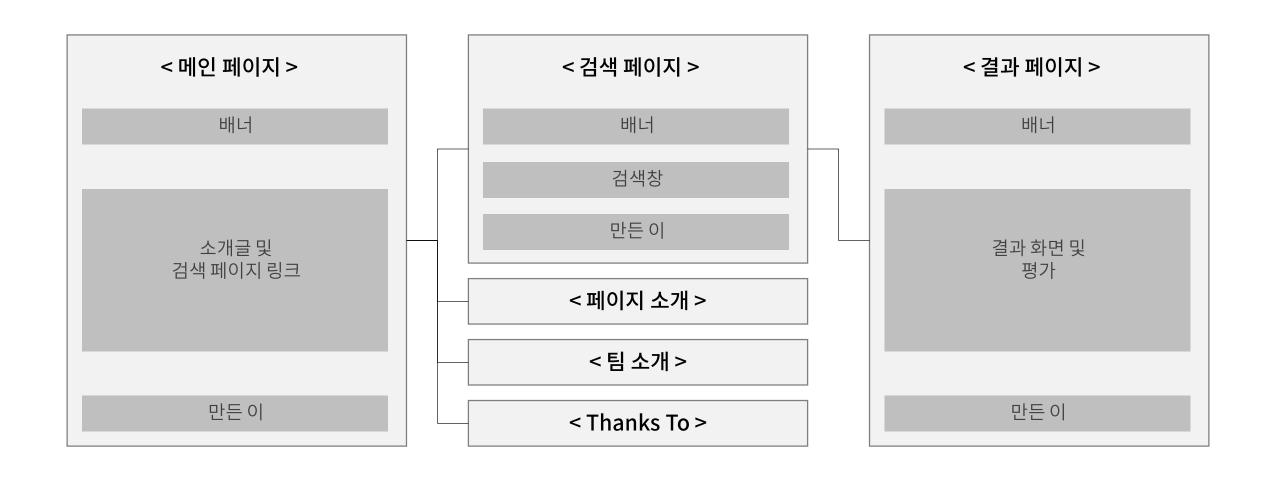


- 일반적인 MVC(Model-View-Controller) 구조
- 파일을 수동으로 만들어야 함
- ORM 기능을 사용하기 위해서는 SQLAlchemy라는 DB 주입 패키지 필요

Flask도 고려대상 었지만 풀스택 프레임워크로 웹 개발 전반을 경험하고 확장성을 고려하여 Django를 선택함

# 02-4. 웹 개발

### 이곡어때 음악 추천 시스템을 웹 구현할 때 고려한 기본 구조는 아래와 같음



## 02-4. 웹 개발

### 웹 구현 시 발생한 이슈는 아래 내용으로 해결함

#### Header와 Footer를 중복 작성하지 않기 위한 Javascript

cdn 보안 문제로 Javascript가 적용되지 않아 html만 보임

서버에 올리면 보이기 때문에 큰 문제는 아니지만 크롬 확장 프로그램인 Web Server for Chrome을 설치하면 로컬에서도 정상적으로 확인할 수 있음



### CSS 및 Javascript를 수정해도 HTML에 적용되지 않는 경우 발생



웹 브라우저에서 HTML을 열면, 수정한 파일이 아닌 캐시로 저장된 내용을 불러와 HTML에 최신 수정 내용을 반영하지 못하는 경우 多

Header 안에서 css의 주소를 넣을 때 '?after'을 입력해서 최신 수정 파일을 불러오도록 함

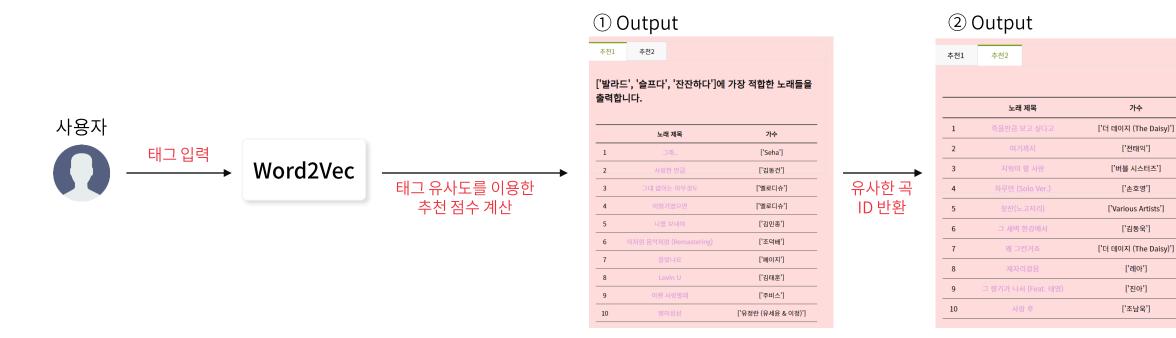
#### 검색 결과 페이지에서 곡 제목과 가수를 표 형식으로 함께 띄우기

- 1) to\_html() 메소드를 이용하여 python 데이터프레임 형식 데이터를 HTML파일로 넘기는 방법을 적용함
  - → 해당 메소드의 style 기본 값이 있어 css 일부가 적용되지 않는 이슈가 있었음
  - → 메소드에 'border=0' 옵션으로 style 기본 값 제거
- 2) HTML 파일로 넘기게 되면 노래 제목에 링크를 다는 것이 불가능해 짐 → to\_html을 사용하는 대신, 데이터프레임에 들어가는 데이터를 zip으로 묶어 넘겨주고 HTML 내에서 for문으로 테이블을 생성함



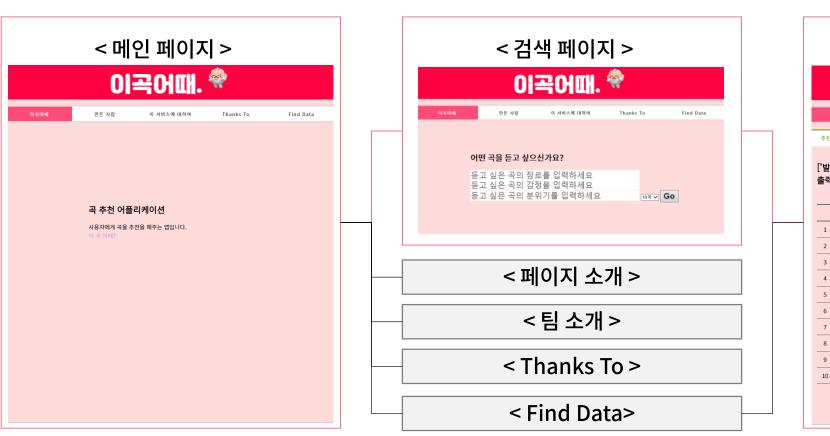
# 03-1. 음악 추천 모델

### 이곡어때 음악 추천 시스템을 실행한 결과는 아래와 같음



# 03-2. 웹 구현

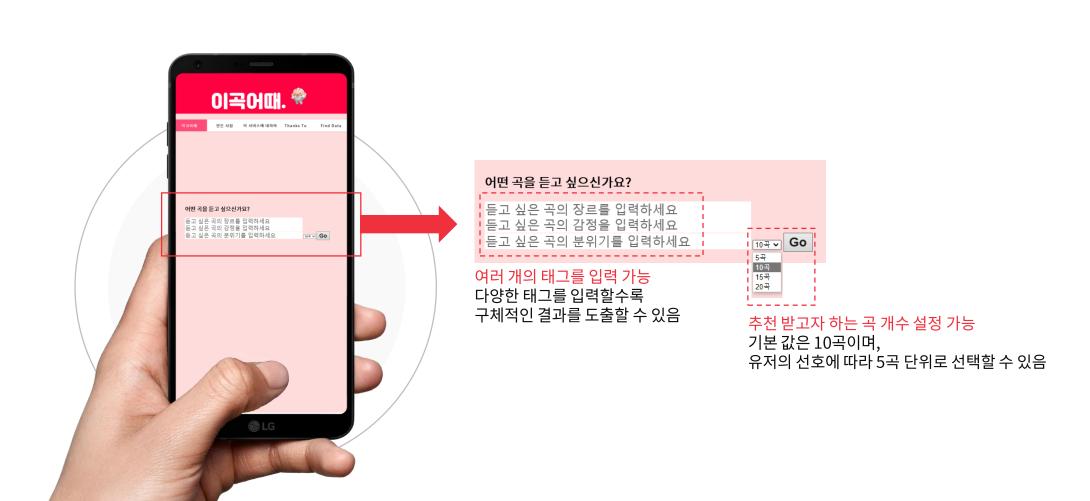
### 이곡어때 음악 추천 사이트는 아래와 같이 구현함





# **03-2. 웹 구현** – 검색 페이지

### 이곡어때 음악 추천 사이트는 아래와 같이 구현함



## **03-2. 웹 구현** – 결과 페이지

### 이곡어때 음악 추천 사이트는 아래와 같이 구현함



# 03-3. 시연 영상



이국어때 만든 사람 이 서비스에 대하여 Thanks To Find Data 계시판

#### 곡 추천 어플리케이션

사용자에게 곡을 추천을 해주는 앱입니다.

이 곡이때?



## 03-4. 한계점

### OOV 문제를 완전히 해결하지 못함

### Word2Vec 모델의 학습 방법으로 인한 OOV(Out of Vocabulary) 문제

- 1) 본 문제를 해결하기 위해 사전 학습된 Word2Vec 모델을 fine-tuning하는 방법을 아래와 같이 적용해 봄
  - 1-1) 한국어 사전학습 모델인 ko.bin을 불러와 fine-tuning 후 모델 적용 시도
    - -> 사전학습 모델 말뭉치에 있는 단어가 결과로 반환되어 노래를 추천할 수 없었기 때문에 <u>이곡어때 시스템을 위한 모델로 적절하지 않다고 판단</u>함
  - 1-2) 사전학습 모델을 fine-tuning한 후 벡터 값만 가져오고, 입력 값으로 들어올 태그들 사이의 코사인 유사도를 계산하면 그 중 유사도가 높은 태그를 관련 태그로 설정할 수 있을 것이라고 생각하여 적용 시도
    - -> 가능했지만 우리 데이터에 없는 단어들로 인해 태그 간 유사도에 왜곡이 생겨 사전학습 모델을 사용하지 않기로 결정함
- 2) 검색 페이지에서 카테고리를 제공함으로써 사용자가 입력할 수 있는 값을 제한하여 오류 발생 가능성을 줄이려는 시도를 했으나, 보다 근본적인 해결방법을 고안해 낼 필요가 있음